

Searching Across the International Space Station Databases

David A. Maluf, Ph.D.
David.A.Maluf@nasa.gov
William J. McDermott
National Aeronautics and Space
Administration
Code TI, Mail Stop 269-4
Moffett Field, CA 94035-1000

Ernest E. Smith
ernest.e.smith@nasa.gov
National Aeronautics and Space
Administration
Code DA4, 2101 NASA Road 1
Houston, TX 77058

David G. Bell Ph.D.
dbell@riacs.edu
Mohana Gurram
Research Institute for Advanced
Computer Science
Moffett Field, CA 94035-1000

Abstract— Data access in the enterprise generally requires us to combine data from different sources and different formats. It is advantageous thus to focus on the intersection of the knowledge across sources and domains; keeping irrelevant knowledge around only serves to make the integration more unwieldy and more complicated than necessary. A *context* search over multiple domain is proposed in this paper to use context sensitive queries to support disciplined manipulation of domain knowledge resources. The objective of a context search is to provide the capability for interrogating many domain knowledge resources, which are largely semantically disjoint. The search supports formally the tasks of selecting, combining, extending, specializing, and modifying components from a diverse set of domains.

This paper demonstrates a new paradigm in *composition* of information for enterprise applications. In particular, it discusses an approach to achieving data integration across multiple sources, in a manner that does not require heavy investment in database and middleware maintenance. This *lean* approach to integration leads to cost-effectiveness and scalability of data integration with an underlying *schema-less* object-relational database management system. This highly scalable, information on demand system framework, called NX-Search, which is an implementation of an information system built on NETMARK. NETMARK is a flexible, high-throughput open database integration framework for managing, storing, and searching unstructured or semi-structured arbitrary XML and HTML used widely at the National Aeronautics Space Administration (NASA) and industry.

Keywords— Intelligent Information Systems, Netmark, NX-Search, Context Search, Semantic Interoperation, Heterogeneous Systems, Integration

TABLE OF CONTENTS

1. INTRODUCTION.....	1
2. BACKGROUND.....	2
3. APPROACH.....	2

4. ARCHITECTURE	3
5. IMPLEMENTATION AND CASE STUDY	4
REFERENCES	6

1. INTRODUCTION

A number of technologies have been developed to support large-scale interoperation among distributed information sources and applications such as databases and heterogeneous sources and applications on the Internet. However, managing large-scale interoperation of *data sources* in an enterprise remains a task, which requires many levels of expertise and an adherence to standards. Many existing information systems and applications have strong notions of interfaces. Efforts like the eXtensible Markup Language, the next generation of HTML, are only ways to put structure into a document, a spreadsheet or a database. These added interfaces allow the specification of the domain's knowledge or the domain component syntax with the data.

Leading efforts in interoperating among multiple domains are often implemented as the *union* of multiple domains. An immediate increase in data integration is noticeable to data integration architects. Yet, there are negative economical consequences to the maintenance of the union of these interfaces. The union of multiple domains always turns out to be a low hanging fruit in the beginning, but always turns out not to be efficient for the evolving needs of customers. Worst, however, the union will become impossible to maintain. System evolution, the lifecycle from cradle to grave, now starts to become the most dominant question in the any investment of existing or new information systems.

This paper presents a solution to a problem of interoperability; the illustration for a benchmark is how to interoperate a Microsoft (MS) Word document with spreadsheet, with PDF document and a payroll relational database. Notably, for interoperation and intelligent access to heterogeneous information, the focus should be on the *intersection* of the knowledge, since intersection will define the required articulations. The term articulation refers to the linkages, which join concepts across domains.

The emergent need to define articulations between data resources has been in demand at NASA. We extend and generalize the identification of the articulation to a set of manipulations, such as selecting, combining, extending, specializing, and modifying components from diverse, common and domain-specific collection of sources. To deal with most semantic issues, a *context* and *content* domain selection is proposed. The *intention* is to support disciplined manipulation of resources. The representation of vocabularies and their structure is termed *domain knowledge* delineating the underlying contexts whereas the operations that combine and partition the domain knowledge in a sound and well-behaved manner are termed a *context algebra*. The basic algebra consists of three operations, namely intersection, union and difference (negation is considered an alternate form of the difference). Knowledge in this paper is limited to the knowledge that an expert can extract from a domain and not the domain itself e.g. complete schema dump. The objective of a context sensitive search is to provide the capability for interrogating many knowledge resources, which are largely semantically disjoint, but where articulations can be established on their perspective context. This paper describes the role of context search among multiple interfaces. It also demonstrates the use of a context algebra, which provides users and system developers with the ability to intelligently manipulate components in real time.

2. BACKGROUND

The development of the mediation model reported in this paper is motivated by the need of interoperability among existing domain-specific representation of knowledge (structure and layout) and their respective formalisms (HTML, XML, Objects etc.). The spirit of this paper is to underline practical aspects of retrieving effectively data with an objective of matching or exceeding the agency current needs. However, what follows is a brief review of what is commonly known in the Knowledge Representation and information integration community [1][2].

The series of knowledge representation formalisms and frameworks starting with KL-One and currently culminating in systems and approaches like semantic-web provide powerful tools and knowledge expressiveness. However, they were intended to interoperate, but their complexity grows with the data. How much has to be added to their infrastructure and semantic-rich capability to achieve knowledge interoperability is still unclear. While knowledge representation is thought of as being a way to resolve integration problems, most knowledge representation formalisms have focused on paradigms, which assume an integrated environment and have been careless about managing the exceptions. In our approach, we focus on these exceptions.

From a research and technical point of view, there have been two recent efforts that open up possibilities for meaningful knowledge interoperation: the development of context logic and knowledge interfaces for sharing. The advance in context logic is the notion of translating encoded knowledge relative to its context and hence relates the knowledge to its domain. Advances in knowledge sharing revolve around translating multiple knowledge from one formalism to multiple formalisms. However, the problem of translating many domains into different representations will create several problems. Semantic inconsistencies will arise from the terms and relationships used from the merged domains. Additional inconsistencies occur when the knowledge-content differs both in semantics and in compositional granularity. In addition, the union of multiple domain knowledge includes irrelevant knowledge and the result will be large, unorganized, and disproportionately costly to process.

The recent formal paradigm in the direction of porting knowledge from one representation language into multiple ones is done by XML XSLT. For example XSLT is a mechanism for translating from one XML scheme and syntax into multiple-representation schemes. However, directly translating entire knowledge to any arbitrary representation leads to irrelevant knowledge and semantic inconsistencies, disproportioned in content.

With the success of Hypertext Markup Language (HTML) and large-scale content distribution of heterogeneous information, industry pushed the technology further with the eXtensible Markup Language (XML). XML is primarily intended to meet the requirements of large-scale Web content providers for industry-specific markup, vendor-neutral data exchange, one-on-one marketing, workflow management, the processing of Web documents by intelligent clients, and most meta-data applications.

3. APPROACH

Interoperation became an industry fact with XML. XML is a system of standards and specifications that describe how software components, as being the domain knowledge, can interoperate across networks, languages and platforms. XML allows for client-server interaction between heterogeneous objects distributed over a wide-area network; XML makes meta-information describing the objects in a system and their interfaces available so that it can access other objects. Any object defined in XML can play simultaneous roles at the client and at the server. To reach effective interoperability with multiple databases with thousands of parameters, and for objects to plug and play, schemas have to be instantly discovered and integrated as part of the query and not hard-coded semantic mapping. The hard-coding paradigm of database integration schema is a major flaw in today's integration approaches.

Linking the semantics across databases manually is a tedious engineering job with little scalability and high cost in maintenance. But maintenance is the intention of the integration providers. Reflecting on the complexity of the number of databases, engineering linkages is a prohibitive in cost and schedule. While it is in scope of a schema designer to craft a schema for a database, it is not justifiable to craft a additional schemas to integrate heterogeneous databases; as the number of data sources is not bound to the designer. What works for ten may not extrapolate and work for the hundreds or thousands.

The idea of combining *composition* and context discovery and binding with *declarative interfaces* is complementary. Declarative interfaces are primarily about specifying component context syntax and structure. *Composition* binds the contexts into a new temporal subset of information through an intersection of the user query with data sources, followed by the union. The discovery mode is a promising outcome dealing with component design, component binding, and component semantics. Although simple in nature, this formalism is powerful enough to scale. A common MS Word document itself is treated as an independent source with same citizenship as full-scale relational database or XML document. Fundamentally, each is source *contexts* are based on their published interfaces.

4. ARCHITECTURE

The *context* and *content* mapping as well as the *context algebra* has been implemented at the National Aeronautics and Space Administration as an effort to solve an Information Technology (IT) challenge as well as a cultural challenge when dealing with the vast amounts of corporate data that NASA builds daily. The product tool suite is named NX-Search with two distinct components: *Composition* and the *Application* and a *Development Interface API*. The implemented system is named NX, a tool suite built on NETMARK schema-less concepts. The original purpose of the XDB system and its Application Programming Interface (API) is to enable NASA information systems to do something that could not do otherwise and to retrieve specific and precise information from within the contents of documents spread across disparate systems. This paper extends NX-Search to the community with the appropriate syntax and use cases. This demonstrates a custom NX-Search Query with a general flavor, the developers' needs, and minimal effort to use and maintain. The expectation however, is a demonstration of NX-Search ability to query unstructured information using standard set of programming patterns and practices. The use at NASA has been to retrieve data repositories based on both context and content, recompose new documents from the results of the queries, and publish information to the users. These demonstration queries combine relevant information from different sources into custom documents.

The immediate benefit from NX-Search has been to enable users and developers to select and integrate contents from proprietary electronic information software systems using a standard interfaces. NX-Search ability to design custom queries resolves the burden on existing systems to share information in a way not designed for originally. The primary audience for this technology has been primarily National Aeronautics and Space Administration (NASA). Nonetheless, industry and other NASA partners have adopted the system in their realm of products and services. This expansion has enabled a much larger community of developers or Help Desk engineers and IT engineers, creating, testing, and, troubleshoot queries that produce custom searches for work groups to meet a new generation of industry products and IT needs, with minimal effort. NETMARK, NX and NX-Search are standard-based application that enforces the World Wide Web Consortium Architecture Domain and Internet Engineering Task Force Standards, including the standards for Relational Databases, WebDAV, XML, XPath, [6][7][9].

The operation of NX-Search WebDAV HTTP API is shown in figure below

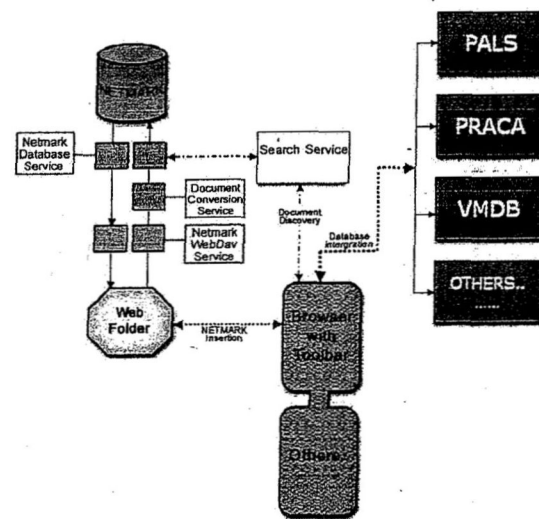


Figure 1 HTTP API – Doc. Insertion, Discovery plus Integrating Multiple ISS Data Bases

The Context+Content Search—takes advantage of the boundaries that demarcate the location of information within a document. NX-Search Context Based Retrieval mechanism is illustrated in figure 2.

A section heading, such as "Procedure", that appears before a paragraph can throw light on the meaning of text within. Aware of the meaning conveyed by section headings. Strong typed databases (semi-structured) uses meta data information as the context for the query, thus enabling context plus content queries.

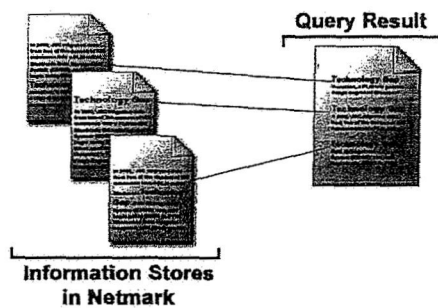


Figure 2 Context-Based Retrieval of relevant text fragments and documents

5. IMPLEMENTATION AND CASE STUDY

International Space Station (ISS) generates large amounts of change information about software that is distributed across multiple documents. This highly distributed information must be securely accessible and usable across the Enterprise. For example ISS Flight controllers needing this information, however, face many challenges. To find the latest information on a pending change to ISS a flight controller must query a range of document types stored in multiple databases, using a different identifier each time. When something doesn't match a flight controller has to sort through multiple of documents of different data type to make the missing connection. Many tools has been developed to aid ISS flight controllers to perform integration and search across many types of documents across multiple sites. And most of these tools has been creating a bigger database or large integrated schema or semantic network. These schemes tended to go obsolete in short period of time as new resources emerged. To address flight controllers' requirements for fast, relevant search capabilities, the approach of dynamic database integration with schema-less have been developed. The Flight Controller leverages this new information storage and retrieval technology to accelerates human intensive searches such as those performed by Station flight controllers. Furthermore, flight controllers require that search results be based on latest documents stored in their original sites a requirement where most centralized and mediated systems failed due to fact of need of a common or routing schema. As a matter of fact, the new system implements a periodical discovery and checking documents at authoritative sites for changes and updating cached documents in provides flight controllers with a level of assurance that information returned is based on latest documents. This new approach tends to eliminate the human engineering bottleneck crafting the needed integrating scheme.

A fold of this paper is to validate novel information integration capabilities in the context of high level of robustness flight environment and work-scenarios such as flight-controllers. These work scenarios dictate the

document types and information sources which are constantly changing and updated. The solution to data currency provides assurance that search results are based on the latest documents from authoritative sources and provide also a long-term solution to data currency requirements. The focus of investigation, and of this paper, is on work scenarios that involve change documents stored in multiple databases. These databases have multiple document-types such word, spreadsheets, PDF, drawings, HTML and standard structured databases.

A typical Scenario

The Station Joint Software Review Board (JSRB) is responsible for managing ISS software changes to ISS. The JSRB uses a Software Change Request (SCR) form, and a number of supporting documents, to track software changes. The SCR and support documents contain inter-connecting links to one another. This rich network of inter-connected documents about a change is logically connected with a major ISS technical category. ISS core areas and other categories make up the structure of most information about Space Station. Station also uses this structure to organize ISS flight-controller groups. A flight controller is trained in one or more technical areas. Later, the flight controller is assigned to a group that is focused on a technical category in ISS structure. From the perspective of software modifications and work-around, flight controllers are primarily concerned with changes pertaining systems on ISS. Consequently, the paper focuses on scenarios that addressed change documentation about this technical area.

The scenario workflow of a documentation of a problem usually starts when ISS Engineering, NASA and its contractor documents a problem in an Item for Investigate (IFI) and sends the IFI to the JSRB. The JSRB determines whether the item in question is a genuine problem that needs to be fixed. If so, the JSRB uses information in the IFI as the basis of a new Software Change Request (SCR). Although the SCR is a central document for tracking problem-resolution and deployment of software changes, the SCR and supporting documents each contain links to one another, defining and connecting relationships between documents. Without some means of preserving relationships within and between documents, it is difficult and time consuming to glean all of the relevant information on an item in question. The number of external sources is not fixed and could spans many other database like problem reporting databases (PRACA), parts databases. Drawing databases, procedure databases etc.

Although the SCR is a central document for tracking problem-resolution and deployment of software changes, the SCR and supporting documents each contain links to one another, defining and connecting relationships between documents. Flight controllers have a particular interested in retrieving information that is distributed across a range of documents about a software change. For example, work-around, release-schedules and documentation changes about


```

graph TD
    RICE[RICE Web Site and Data Base]
    SCR[Software Change Request (SCR)]
    SPS[System Program Store (SPS)]
    SJF[Sched. Job Change Form (SJF)]
    CN1[Change Notice (We for PPL)]
    CN2[Change Notice (CN)]
    MMSDN[MMSDN]
    PRN[Status Printed Release Notice (PRN)]

    RICE --- SCR
    RICE --- SPS
    RICE --- SJF
    RICE --- CN1
    RICE --- CN2
    SJF --- MMSDN
    SJF --- PRN
    
```

SPS is : server/running archive of a time

Without some means of preserving relationships within and between documents, it is difficult and time consuming to glean all of the relevant information on an item in question.

Software Change Requests and related change documents are accessible to flight controllers. Flight controllers who seek to obtain the entire body of relevant information about changes, however, face challenges. Change documentation is stored in a database that provides a limited search capability. Flight controllers can use search tool to query within a document-type: search for SCRs for example with keywords in TITLE attribute. A new query is needed for each source. It is very important to underline that the linkages among source is knowledge of the flight controller experience and not a simply a semantic map as it may seems. Search results of this title do not capture the relationship between the SRC that was found and it's related, supporting documents. Thus, flight controllers can not search across document-types in a single query.

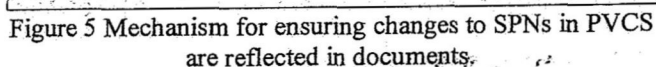
A scenario based on a search across multiple sources and documents is a requirement. The multiple source/document-types in this scenario on are listed in table 1:

Table 1 Sample of Document-types included in investigation

NX-Search fine-grain search-capability when querying change documents and specifications and also other notes, specifications, and reports from a range of information sources. Typically this underlines an integration of documents from multiple information sources matching different granularities.

The IFI database, CHITS, and DF's Anomaly Report contain information related to changes. To search across the range of document types stored in multiple systems and these systems requires software that is sensitive to relationships of sections within documents as well as relationships across documents.

The high throughput of NX-Search back-end became also the solution to the data currency problem. This is done by synchronizing updates to originals with documents cached. A periodical pedigreed cache becomes thus available. A logical representation of the update process is shown in figure 5.



5

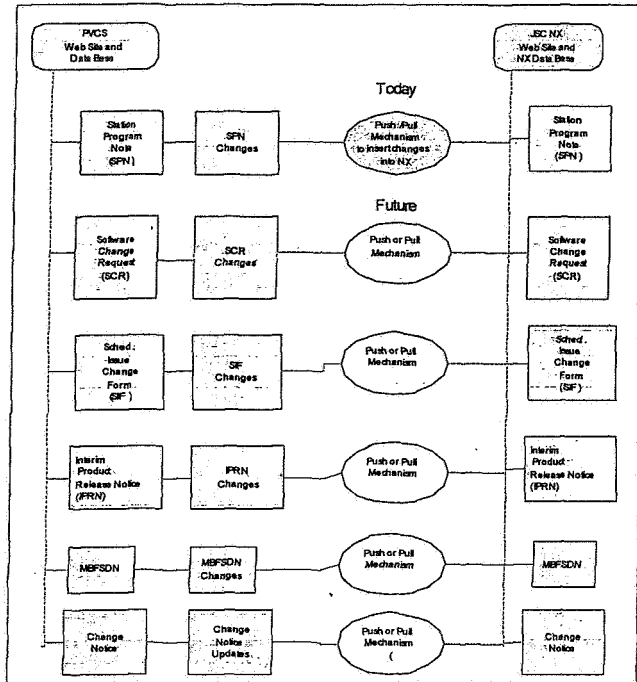


Figure 6 Generalizing the SPN update mechanism to other documents

The cache captures changes to SPNs and stores the information in NX-Search. The update reflect changes in these sources monitoring the evolution of the remote remote resources. However, the update routine is based on the same mechanism employed by ISS to update information on SPNs stored in the system of records to gain pedigree sign-off for the cached data. Systematically, search results based on the integrated view have the same fidelity to documents in as in their original sources interfaces. This degree of fidelity is highly important to the success of an information system when dealing with flight systems.. Producing a testable capability to for synchronizing changes in multiple sources and NX-Search provides a firm line of evidence that the data currency problem is a tractable one. Producing a round about is an essential component in information integration which are asynchronous in nature.

REFERENCES

- [1] A. Halevy et al., Enterprise Information Integration: Successes, Challenges and Controversies ACM SIGMOD International Conference on Management of Data, Baltimore MD 2005.
- [2] D. Maluf, P. Tran, NETMARK: A Schema-Less Extension for Relational Databases for Managing Semi-structured Data Dynamically ISMIS 2003.
- [3] J. Shanmugasundaram et al., SIGMOD Record 30, 20-26 (2001).
- [4] H.V.Jagadish et al., VLDB Journal 11, 274-291 (2002).
- [5] Oracle 9i Database Release 9.0.1 Developer Guide.
- [6] J. Whitehead, Y. Goland, WebDAV: A network protocol for remote collaborative authoring on the Web CSCW 1999.
- [7] XSLT, <http://www.w3.org/TR/xslt>.
- [8] D. Maluf, P. Tran, T. La, "An Extensible 'Schema-less' Database Framework for Managing High-Throughput Semi-Structured Documents IASTED, Applied Informatics 2003.
- [9] Xalan, <http://xml.apache.org/xalan-j/>.
- [10] D. Draper, A. Y. Halevy, D. S. Weld., The Nimble XML Data Integration System ICDE 2001.
- [11] C. A. Knoblock et al., International Journal of Cooperative Information Systems (IJCIS) Special Issue on Intelligent Information Agents: Theory and Applications 10, 145-169 (2001).